

# Сколько стоит разработать собственную сборку Hadoop

История и техника как получилось в Сбере

Семён Орлов



# Привет, я Семён Орлов



СБЕР



Product Owner SDP Hadoop



2 ВО: физика полупроводников  
и менеджмент



1/3 жизни в ИТ.  
До прихода в Сбер работал в



1/6 жизни посвятил работе с

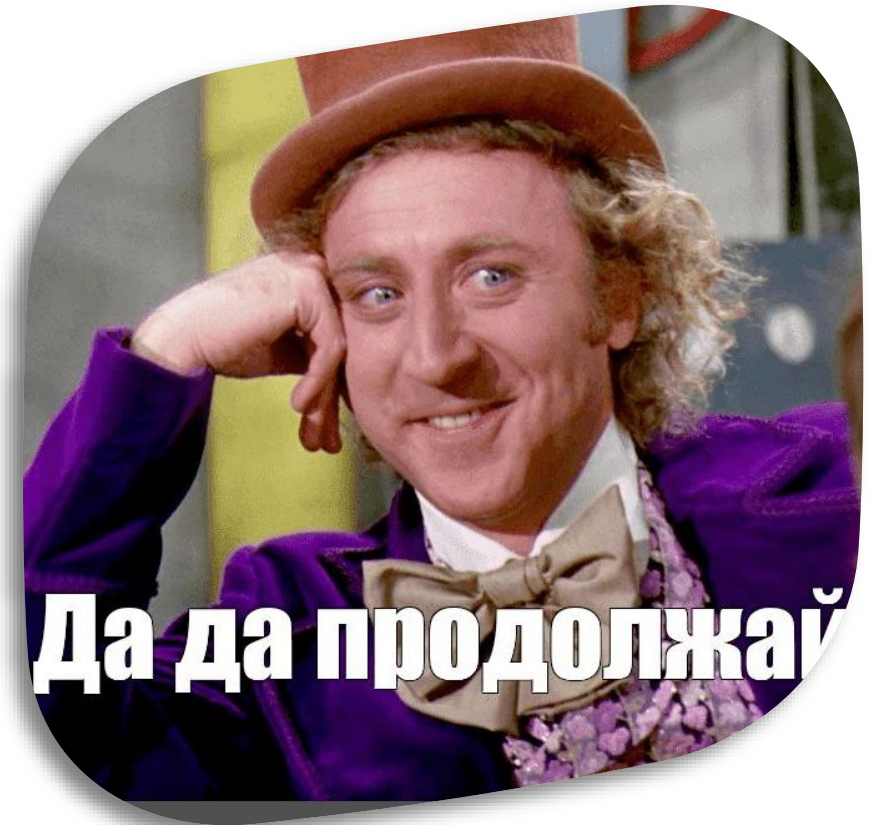


`hdfs dfs -setrep -w 3`

# О чём расскажу?

☆ Историю, как делали

Если нужна выжимка, мотай +25 минут



# Масштабы сегодня



БАШНЯ ФЕДЕРАЦИЯ



**150** ПБ  
ДАННЫХ

**27** ПБ  
В КЛАСТЕРЕ

**155**  
КЛАСТЕРОВ

**374** УЗЛА  
В КЛАСТЕРЕ



# Окружение 2018 года

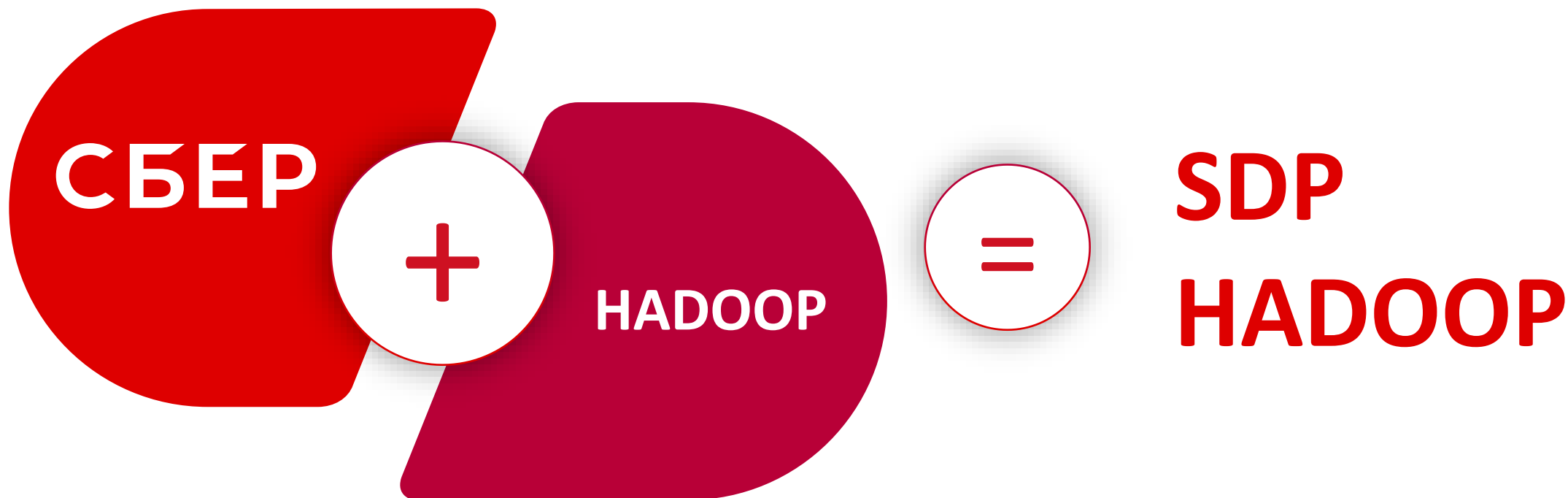


**1 trillion \$  
Company**

# Окружение



# Окружение лифты



# Зачем начали делать? Мотивы

## Поддержка

- Сложная схема поддержки
- Невозможность работы on site
- Затраты на архитектурную поддержку

## ТСО

- Существенная доля - лицензии
- Снижение стоимости 1ТБ

## Кастом

- Безопасность
- Интеграции в КАП
- Автоматизация

## Экосистема

- Sber#
- Cloud
- Гостех



# Продать идею бизнесу

1

ТСО

2

Поддержка

3

Плюшки

- ✓ Скорость обновлений
- ✓ Компоненты
- ✓ Вендор-риск

# Получить ресурсы

# Ключ: команда для «узких» кейсов

2000

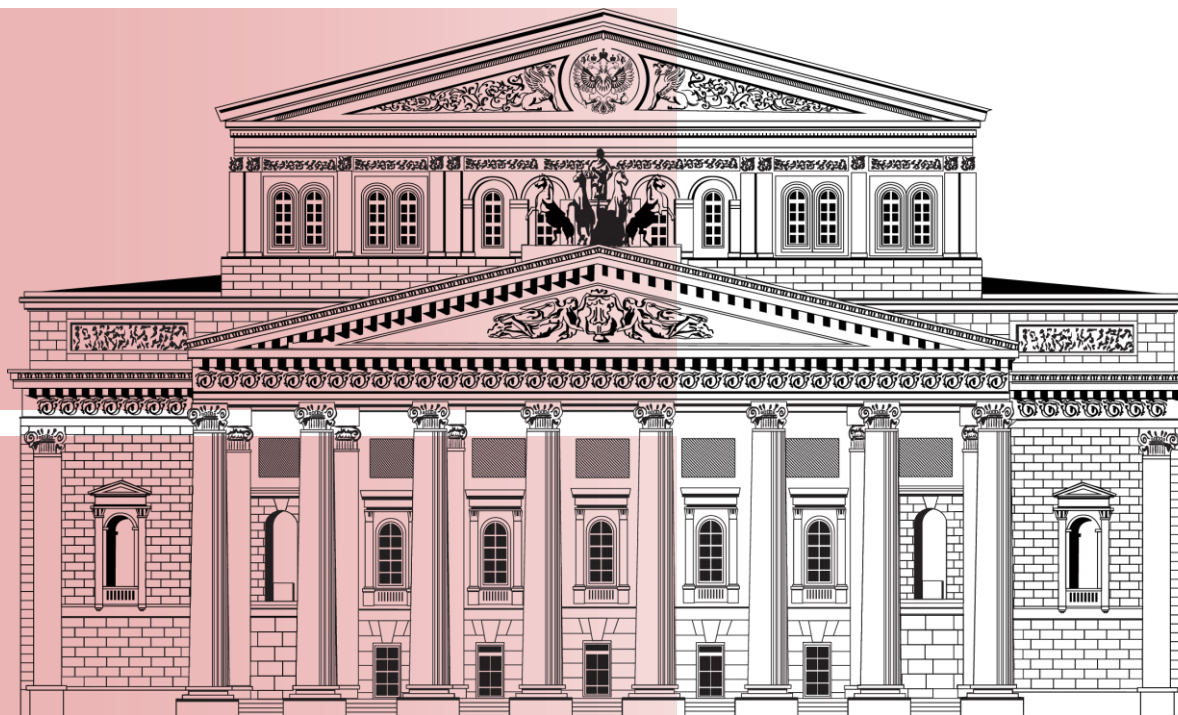
x2

**879**

зрительских  
мест

$\pm 1\ 200 + 500 +$   
 $+ 300$  vs 10

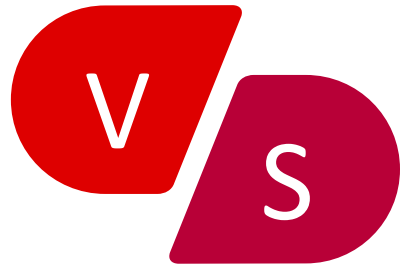
500/36 vs 10/15



# Метрики или результаты: что важнее для выживания

## Базовые метрики

- MAU/DAU
- CPI
- ARPU
- LTV
- RR



## Результаты

- Кластеры в PROD
- Бесшовная миграция
- Администраторы могут сопровождать
- Безопасность не дырявая

# EOL: **product vision** не поможет



Идея



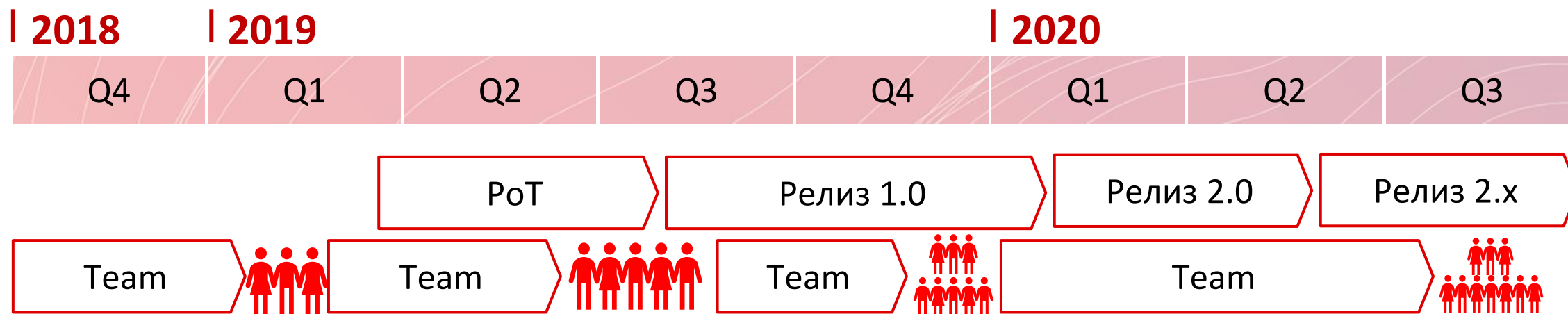
Оно работает



EOL Hadoop  
2



PROD  
кластер



# Собрать мало, **мигрируй**

| 2019

| 2020

| 2021

| 2022

H1

H2

H1

H2

H1

H2

H1

H2

PoT

Релиз 1.0

Релиз  
2.0

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

Релиз  
2.x

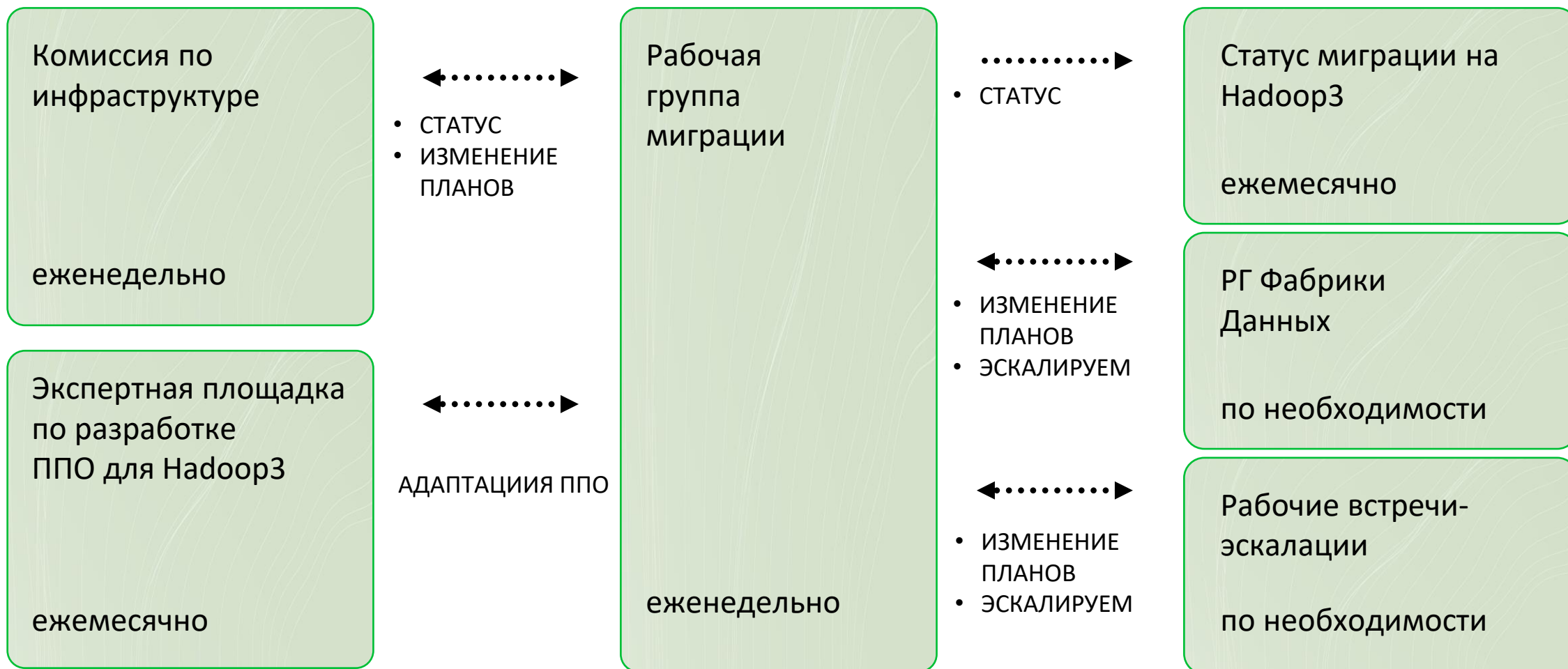
Prepare

Миграция vendor-> SDP



# Проект миграции

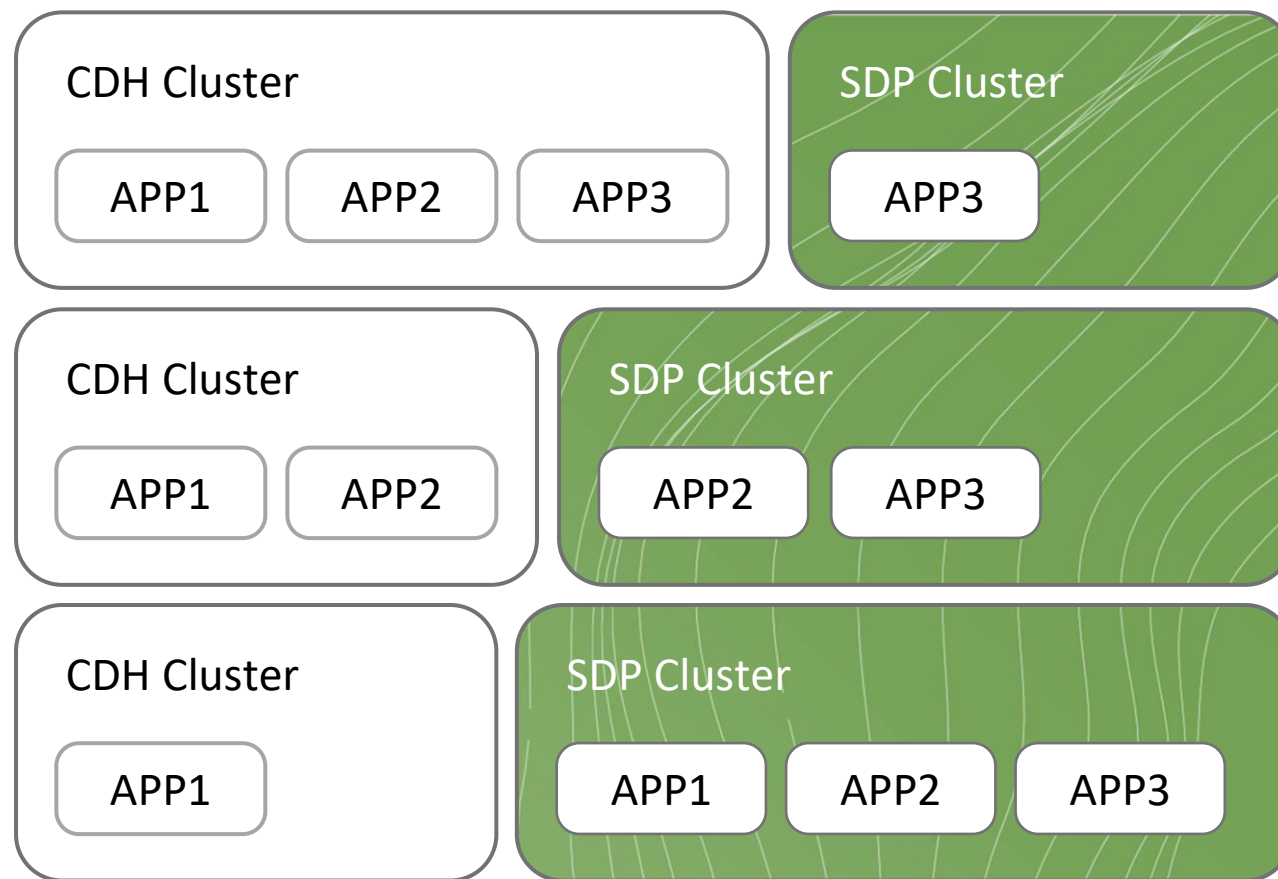
100+ кластеров | 1 000+ экземпляров ППО | 50 тех сервисов



# Схема **миграции**

ПЯТНАШКИ  
С СЕРВЕРАМИ

**+20%**



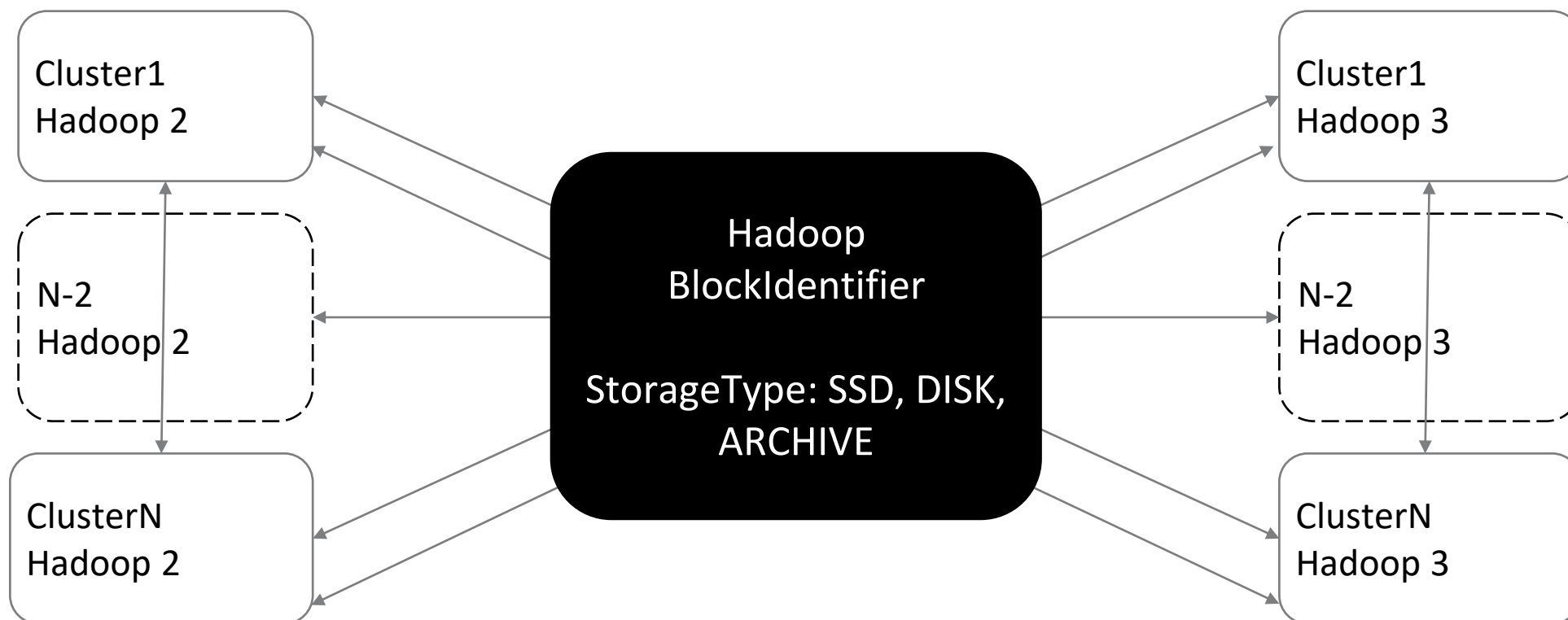
# Какие задачи на Hadoop

А что дорабатывали?

- DWH
- Реплики batch и NRT
- 9 фреймворков
- DataMarts
- ML
- Integrations: a lot



# Сложности

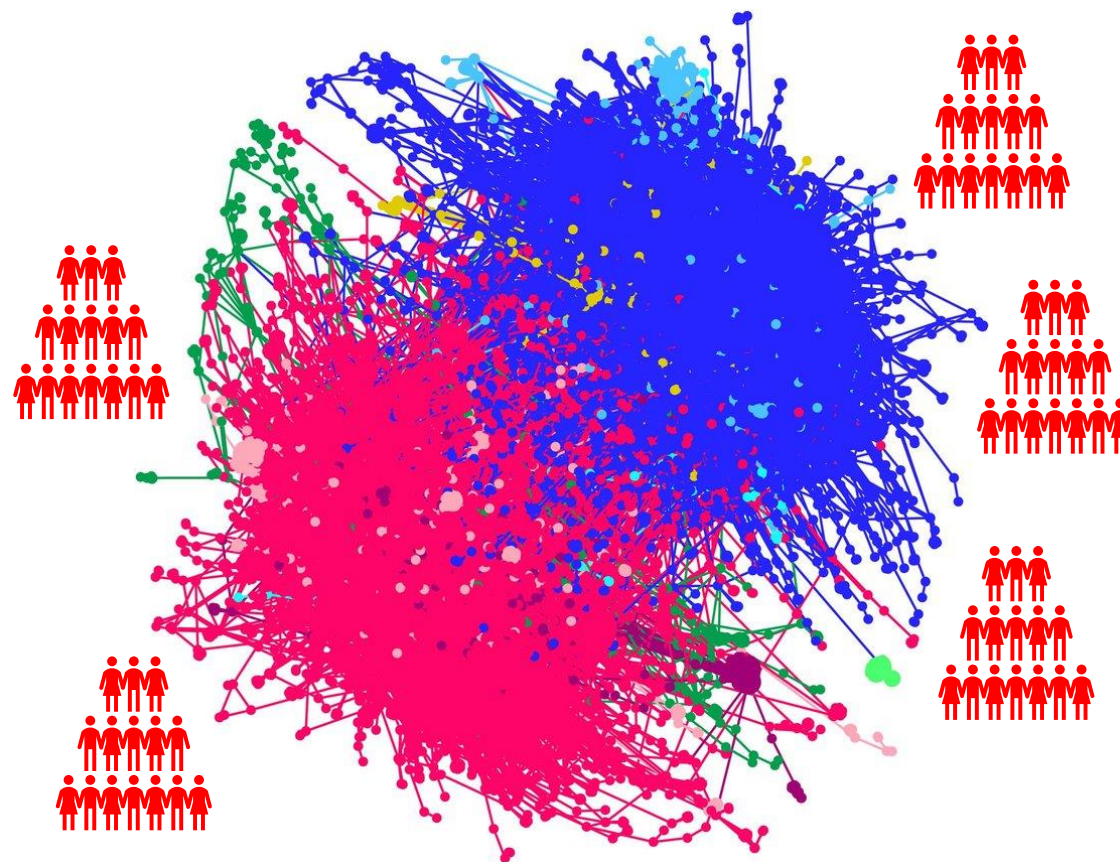


# Федерации 100х кластеров

Information Lifecycle  
Management

Proxy, Federation

JBOD  → RAID 1 





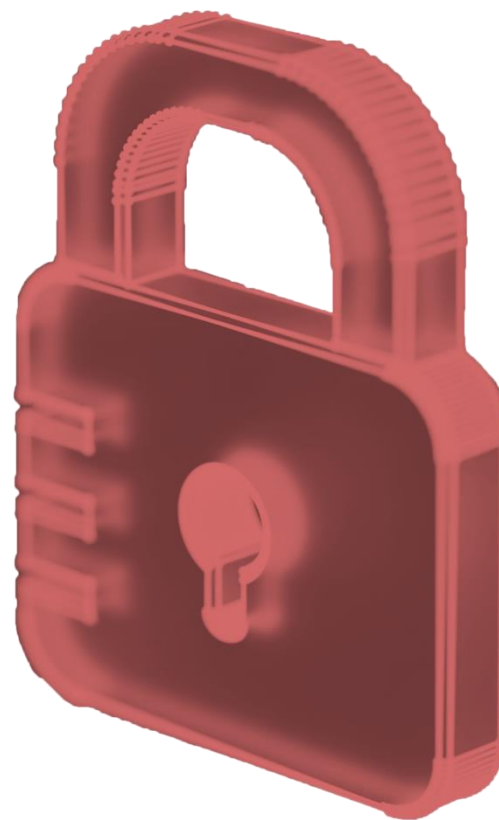
# Доработки кибер-безы

Грануляция привилегий hdfs snapshot

CVE-фикс для Spark 3

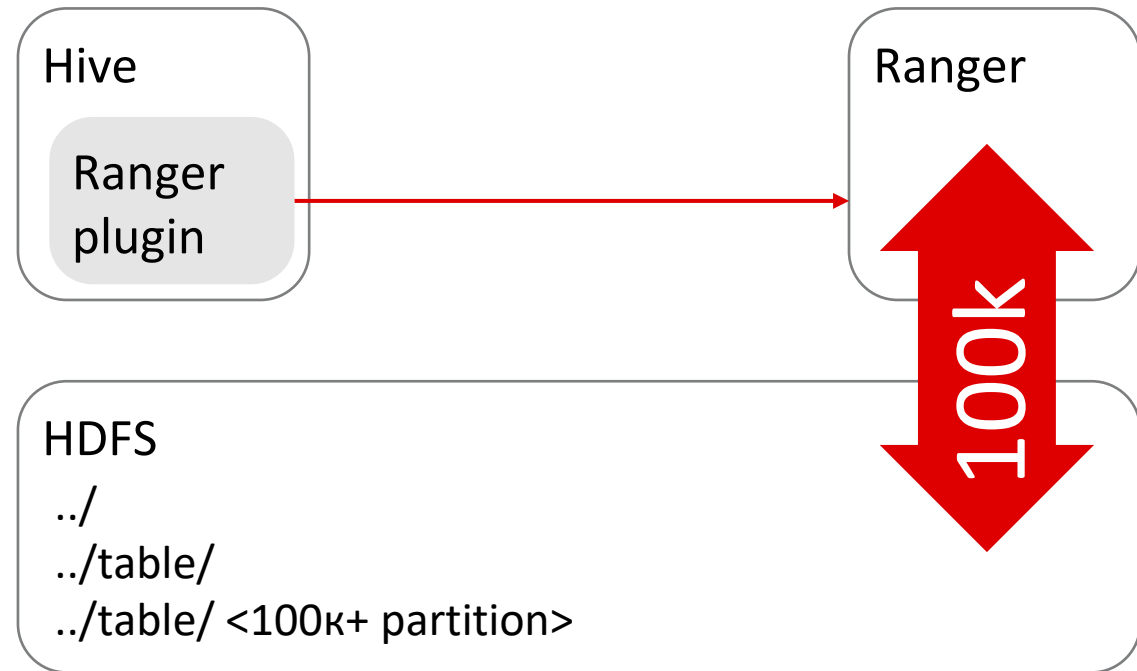
Аудит Ambari

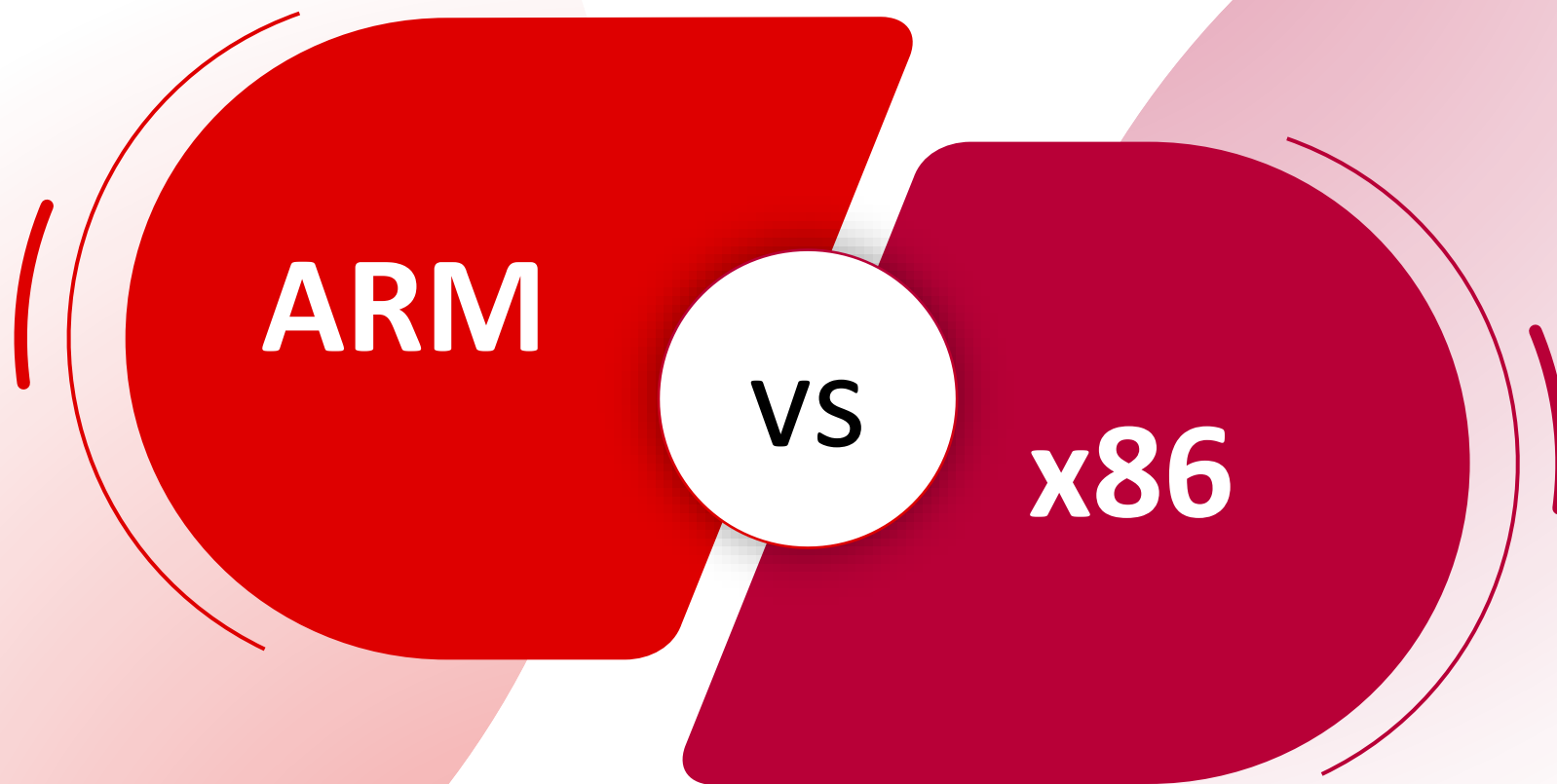
Совместимость Ranger и Hbase



# Доработки Hive

create external table  
big\_table





# Как определить что пора?

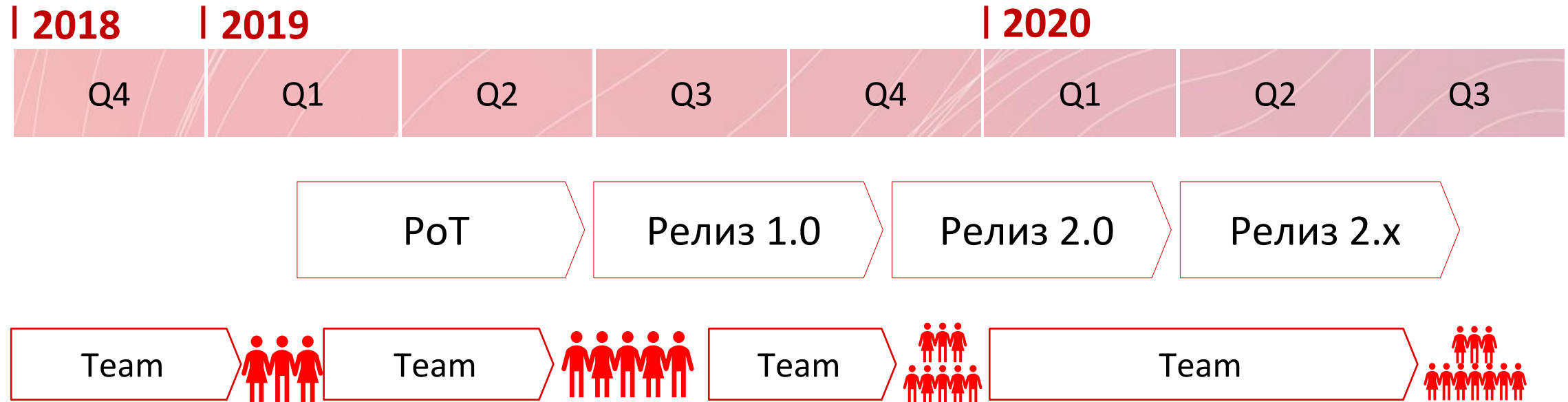
Потребности

Возможности

Сроки

Cost 2 value-анализ

# Если бы пришлось повторить



$$24 - 22 - 3 - 34 - 4 = 12$$



# Какую команду **собирать**



MVP



DevOps



QA



LAUNCH



DevOps



QA



DM



PROD



DevOps



QA



DM



DEV



DE

# Выводы

1

Верить  
в идею

2

Если делать, то  
получится

3

Есть **Команда** -  
есть результат

# СПАСИБО

[t.me/samorlov](https://t.me/samorlov)



# Сцена после титров

1. Что будет с классическим дистрибутивом ближайшие 3 года
2. Как трансформируется платформа обработки данных в Сбере
3. Когда Hadoop действительно умрет или переродится